# Title:
Subtitle

©Patrick Hall 20XX

May 15, 2024

**Abstract**

## 1 Introduction

The National Institute of Standards and Technology Artificial Intelligencde (AI) Risk Management Framework (RMF).[1]

## 2 Generative AI Governance

## 3 Generative AI Inventories

## 4 Generative AI Risk Tiers

## 5 Generative AI Risk Measurement

## 6 Generative AI Risk Management

## Conclusion

## Acknowledgments

## Abbreviations

- AI: Artificial Intelligence

- AI RMF: Artificial Intelligence Risk Management Framework

- GAI: Generative AI

- RMF: Risk Management Framework

[1] NIST AI. Artificial Intelligence Risk Management Framework (AI RMF 1.0). 2023.